



# Community mining with graph filters for correlation matrices

Pierre Borgnat, Paulo Gonçalves, Nicolas Tremblay, Nathanaël Willaime-Angonin

## ► To cite this version:

Pierre Borgnat, Paulo Gonçalves, Nicolas Tremblay, Nathanaël Willaime-Angonin. Community mining with graph filters for correlation matrices. Asilomar Conference on Signals, Systems, and Computers, Nov 2015, Monterey (CA), United States. hal-01245926v2

**HAL Id: hal-01245926**

**<https://inria.hal.science/hal-01245926v2>**

Submitted on 18 Dec 2015

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# COMMUNITY MINING WITH GRAPH FILTERS FOR CORRELATION MATRICES

Pierre Borgnat<sup>(1)</sup>, Paulo Gonçalves<sup>(2)</sup>, Nicolas Tremblay<sup>(3,4)</sup>, Nathanaël Willaime-Angonin<sup>(1,5)</sup>

<sup>1</sup> CNRS, Laboratoire de Physique, École Normale Supérieure de Lyon, France

<sup>2</sup> INRIA, LIP, DANTE team, École Normale Supérieure de Lyon, France

<sup>3</sup> INRIA Rennes - Bretagne Atlantique, Beaulieu Campus, Rennes, France

<sup>4</sup> Institute of Electrical Engineering, EPFL, Lausanne, Switzerland

<sup>5</sup> École Normale Supérieure (Student), Paris, France

Corresponding author: pierre.borgnat@ens-lyon.fr

## ABSTRACT

Communities are an important type of structure in networks. Graph filters, such as wavelet filterbanks, have been used to detect such communities as groups of nodes more densely connected together than with the outsiders. When dealing with times series, it is possible to build a relational network based on the correlation matrix. However, in such a network, weights assigned to each edge have different properties than those of usual adjacency matrices. As a result, classical community detection methods based on modularity optimization are not consistent and the modularity needs to be redefined to take into account the structure of the correlation from random matrix theory. Here, we address how to detect communities from correlation matrices, by filtering global modes and random parts using properties that are specific to the distribution of correlation eigenvalues. Based on a Louvain approach, an algorithm to detect multiscale communities is also developed, which yields a weighted hierarchy of communities. The implementation of the method using graph filters is also discussed.

**Index Terms**— community detection, modularity, correlation matrix, hierarchical communities, graph filters

## 1. INTRODUCTION

A first goal when analysing a set of related signals, such as times series acquired by sensor networks (see Fig. 1), economical series, or any dynamical quantity measured at different points in space, is to discover relations between them, and to group together similar series, before further processing. We adopt a network model to study this question, assuming that each series represents a node of the network and that the relation between any two series is the weight of the corresponding edge in the network. Then, groups of closely related series will appear as communities in this network, i.e., groups of nodes having a larger proportion of links within the group than without [1]. Existence of communities is a frequent and well studied feature of complex networks [2].

The objective of the present work is to show that one can cluster together series, even if they are correlated and nonstationary, considering the correlation matrix of the whole collection of series, importing thus, the concept of communities from network analysis. However, as it was shown in [3], correlation matrices are not adjacency

matrices of networks and the classical modularity metrics [2] has to be adapted. Recalling how classical modularity is extended to correlation matrices, we first show some resulting pitfalls in community detection. More precisely, we will be confronted to two common problems of community detection: one is the presence of global modes or trends among all the individuals (the so-called market mode in economy) that can mask the specific relations within groups. A second difficulty comes from the size heterogeneity of the groups which raises the tricky question of resolution limit for modularity in usual networks [4]. In a second step, we propose solutions to avoid such pitfalls. This leads us to a new algorithm for multiresolution community mining from correlation matrices. The proposed approach is tested both on simulated examples and on a real-world example of temperature sensor networks.

## 2. BACKGROUND

### 2.1. Community detection with modularity matrix

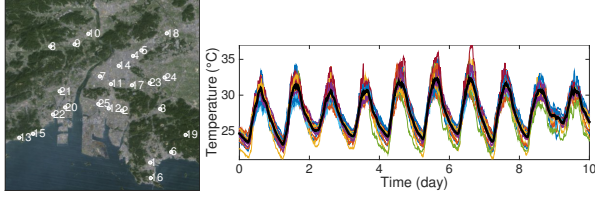
Numerous works and methods exist to find communities in complex networks, many of which being reviewed in the survey of S. Fortunato [1]. As a loose definition, a community is a set of nodes that has a larger number of links inside the group than with the outside. That said, it remains to decide on a precise metrics to quantify this property. Methods have been proposed ranging from the use of spectral clustering or cut algorithms (see the review in [5]) to the use of the popular modularity metrics [2], information-theoretical approaches [6], or graph wavelet based methods where one relies on graph filterbanks defining wavelets to provide ego-centered views from each node [7].

The starting point here is **the modularity** of a network. This quantity measures how relevant a given node partition is to represent the different communities that compose the network. It is calculated by comparing the strength of the edges within communities to a null model corresponding to a random rewiring of the links while the nodes degree is kept unchanged. For a partition described by  $\sigma_i$  (the label of the group of node  $i$ ), the modularity is defined as

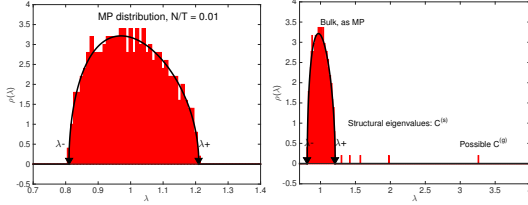
$$Q(\sigma) = \frac{1}{2m} \sum_{ij} (\mathbf{A}_{ij} - \langle \mathbf{A}_{ij} \rangle) \delta(\sigma_i, \sigma_j) \quad (1)$$

where  $\delta$  stands for the Dirac function,  $\langle \mathbf{A}_{ij} \rangle = \frac{k_i k_j}{2m}$  with  $k_i = \sum_j \mathbf{A}_{ij}$ , and  $2m = \sum_i k_i$ . A good partition in communities is then associated to a large value of  $Q(\sigma)$ . More concisely, the modu-

This work was partly funded by the European Research Council, PLEASE project (ERC-StG-2011-277906); and the ANR-14-CE27-0001 GRAPHSSIP grant.



**Fig. 1.** Live E! sensor network [11] in Kurashiki city, Okayama Prefecture, Japan, and the temperatures over 25 sensors for 10 days.



**Fig. 2.** Estimation by histograms of the eigenvalue distribution of  $\mathbf{C}$  for  $N/T = 0.01$ , (left) when following the hypotheses leading to the Marchenko-Pastur distribution (also shown in black above) and (right) when there is a structured mode with eigenvalues outside the bulk (and possibly a global mode as well).

larity (1) can be rewritten as follows:

$$Q(\sigma) = \frac{1}{2m} \text{Trace}(\sigma^T (\mathbf{A} - \langle \mathbf{A} \rangle) \sigma) \quad (2)$$

where  $\sigma$  is a matrix coding for the communities,  $\sigma_{ij} = 1$  if node  $i$  is in community  $j$ , and 0 otherwise.

Maximization of  $Q(\sigma)$  is hard as it requires an optimization procedure over the huge space of all partitions of any size of the network. There are several possibilities to find the partitions that (approximatively) maximize this modularity: simulated annealing, spectral methods [1, 2]. The modularity matrix  $(\mathbf{A} - \langle \mathbf{A} \rangle)/2m$  can be studied in itself to find a relevant partition in a community (see [8, 9]). In the following, we will use the sub-optimal yet very efficient method that is the greedy Louvain algorithm [10].

This **Louvain algorithm** will be used as it opens the way to the study of large scale problems [10]. In a nutshell, this algorithm iterates two steps:

- (1) Select a node and group it with the node that causes the largest increase of  $Q$ ; do this sequentially with all other nodes.
- (2) Merge the nodes of the same community to form a new network whose nodes are the communities formed at step (1).

At the beginning of each step (1), the partition assigns one node to one community (singletons), then the different communities grow or disappear by absorption. The algorithmic efficiency comes from the fact that it is possible to write the modularity increase due to steps (1), without considering the entire graph but only the nodes to be grouped together, and that it is possible to consistently and easily derive the weighted edges of the pruned graph from the initial weighted edges.

## 2.2. Decomposition of correlation matrices

As argued in [3], correlation matrices are not adjacency matrices of networks and therefore modularity is not readily applicable. In

particular, they show that the null model  $\langle \mathbf{A}_{ij} \rangle$  is failing for correlation matrices and leads to biases. To overcome the problem, [3] proposes to rely on random matrix theory in order to have a relevant null model for correlation matrices. We briefly recall their approach.

Let us consider the time series  $X_i(t)$  for  $i = 1, \dots, N$  at times  $t = 1, \dots, T$ . Their correlation matrix  $\mathbf{C}$  is:

$$\mathbf{C}_{ij} = \text{Corr}(X_i, X_j) = \frac{X_i X_j^T}{\sqrt{\text{Var}(X_i) \text{Var}(X_j)}} = \tilde{X}_i \tilde{X}_j^T \quad (3)$$

with centered and normalized series  $\tilde{X}_i$ . Results from random matrix theory (see, e.g., [12, 13]) can be applied to decompose the correlation matrices in several parts: a random part (or bulk) associated to some null hypothesis, a structured part which describes non-random behaviors and possibly a global (also called market) mode that represents a general mode common to all series (such as the one apparent on Fig. 1).

The random part is associated to the null model that would coincide with stationary, white (i.e., uncorrelated in time) series that are i.i.d. with 0 mean and identical variance. In that case and in the limit of  $N$  and  $T$  large with ratio  $N/T > 1$  still finite, the eigen-decomposition of the correlation matrix  $\mathbf{C}$ , yields eigenvalues following the famous Marchenko-Pastur (MP) distribution:

$$\rho(\lambda) = \frac{T}{N} \frac{\sqrt{(\lambda_+ - \lambda)(\lambda - \lambda_-)}}{2\pi\lambda} \quad (4)$$

for  $\lambda \in [\lambda_-, \lambda_+]$  where  $\lambda_{\pm} = \left(1 \pm \sqrt{\frac{N}{T}}\right)^2$ , and  $\rho(\lambda) = 0$  outside this interval. Fig. 2 shows this theoretical distribution superimposed to the sample histogram obtained for stationary, white random series. This suggests that the correlation eigenvalues lying within the interval  $[\lambda_-, \lambda_+]$  could correspond to a random null model with independent, stationary and white series. Hence, if there is a correlation structure that can be detected, its footprint is expected on eigenvalues that lie outside this interval.

This argument leads to decompose a correlation matrix  $\mathbf{C}$  in its spectral domain. Let's diagonalize  $\mathbf{C}$  (always possible as it is a definite positive matrix) as:

$$\mathbf{C} = \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_N) \mathbf{U}^T, \quad (5)$$

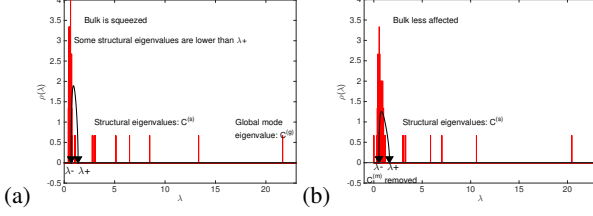
where the eigenvalues  $\lambda_k$  are sorted in decreasing value:  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{N-1} \geq \lambda_N$ . The decomposition comprises three parts:

$$\mathbf{C} = \mathbf{C}^{(r)} + \mathbf{C}^{(s)} = \mathbf{C}^{(r)} + \mathbf{C}^{(s-)} + \mathbf{C}^{(g)}. \quad (6)$$

- $\mathbf{C}^{(s)} = \mathbf{U} \text{diag}(\lambda_1, \dots, \lambda_s, 0, \dots, 0) \mathbf{U}^T$ , where  $\lambda_s$  is the smallest eigenvalue that is still larger than  $\lambda_+$  (there are  $s$  such eigenvalues). It is called the structure mode.
- $\mathbf{C}^{(r)} = \mathbf{U} \text{diag}(0, 0, \dots, 0, \lambda_{s+1}, \dots, \lambda_N) \mathbf{U}^T$ , where  $\lambda_{s+1}$  is the first eigenvalue smaller or equal to  $\lambda_+$ . (there are  $N - s$  such eigenvalues). It is the random mode or “bulk”.
- If there is a so-called global mode due to a general dynamics common to all series, it is associated to the largest structured eigenvalue  $\lambda_1$ . In this case, one can split the structured part in two and remove the global mode  $\mathbf{C}^{(g)}$  from it to keep a (reduced) structure mode:

$$\mathbf{C}^{(g)} = \mathbf{U} \text{diag}(\lambda_1, 0, \dots, 0) \mathbf{U}^T = \lambda_1 \mathbf{U}_1 \mathbf{U}_1^T \quad (7)$$

$$\mathbf{C}^{(s-)} = \mathbf{U} \text{diag}(0, \lambda_2, \dots, \lambda_s, 0, \dots, 0) \mathbf{U}^T. \quad (8)$$



**Fig. 3.** Example of the distribution of the eigenvalues of  $C$ , (a) when there is a global mode and the bulk is squeezed and is not located in the interval  $[\lambda_-, \lambda_+]$  (two isolated eigenvalues that should be in  $C^{(s)}$  are lower than  $\lambda_+$ ), and (b) after filtering of the global mode (the interval  $[\lambda_-, \lambda_+]$  is more representative of the bulk).

### 2.3. Community detection from correlation matrices

Combining results from sections 2.1 and 2.2, the authors of [3] propose to maximize a modified modularity adapted to correlation matrices. If there is no global mode, the modularity is written, mutatis mutandis, as eq. (2) where  $C^{(r)}$  takes the role of the null model:

$$Q_C(\sigma) = \frac{1}{C_{tot}} \text{Trace}(\sigma^T (C - C^{(r)}) \sigma). \quad (9)$$

Moreover, if there is a global mode in  $C^{(g)}$ , they propose to remove it from the correlation matrix and the modularity of  $C^{(s-)}$  reads

$$Q_C^-(\sigma) = \frac{1}{C_{tot}} \text{Trace}(\sigma^T (C - C^{(r)} - C^{(g)}) \sigma). \quad (10)$$

To maximise the modularity  $Q_C(\sigma)$  or  $Q_C^-(\sigma)$ , they resort to the Louvain algorithm described in Section 2.1.

## 3. PITFALLS IN COMMUNITY DETECTION FROM CORRELATION MATRICES

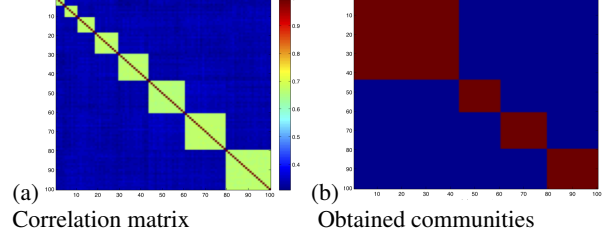
Let us now stress some limits of this approach, using for that, time series  $X_i(t)$  that follow the same model as in [3]:

$$X_i(t) = a \alpha(t) + b_i \beta_{\sigma_i}(t) + c \gamma_i(t) \quad (11)$$

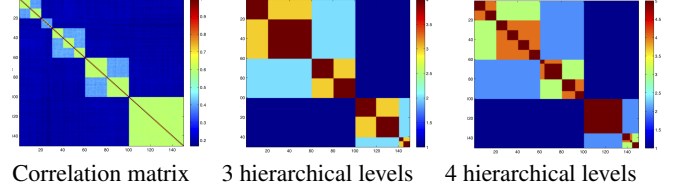
where  $\alpha(t)$  is the global mode (with amplitude  $a$ ),  $\beta_{\sigma_i}(t)$  is the discriminant mode of the group  $\sigma_i$  and common to all time series therein (with amplitude  $b_i$ ), and  $\gamma_i(t)$  is the noise for the node  $i$  (with amplitude  $c$ ). All modes are i.i.d. centered, normalised, white gaussian noises. The following sections describe situations where the proposed method fails at finding the true communities of the model.

### 3.1. Presence of a strong global mode

When there is a strong global mode, the bulk is squeezed as compared with the expected MP distribution and this leads both to misplace  $\lambda_+$  and to have the eigenvector associated to the largest eigenvalue to dominate in  $C^{(s)}$ . See illustration of Fig 3 (a). As a consequence, the Louvain algorithm usually finds only one community: all series share the same global evolution (e.g. see series of Fig. 1). If we were to know that there is a global mode, we could try to update the (wrong) value of  $\lambda_+$ . For instance, considering the one-class model  $X_i(t) = a \alpha(t) + c \gamma_i(t)$  leads to a bulk ending in  $c\lambda_+/\sqrt{a^2 + c^2}$  instead on  $\lambda_+$ . However, the parameters  $a$  and  $c$  are unknown beforehand in practical situations and this is not easily possible.



**Fig. 4.** Resolution limit of modularity: given the model correlation matrix on the left, the maximization of modularity outputs the communities on the right, merging together smaller groups.



**Fig. 5.** Hierarchical mining of communities, commented in 3.2.

### 3.2. Limit in resolution

A more standard weakness of modularity maximisation is its inability at finding small scale structures.[4]. This is illustrated in Fig 4 where  $N = 100$  series are grouped in communities of sizes: 4, 6, 8, 11, 14, 17, 19, and 21. The community detection gathers the first 5 groups into a unique global community, because of the natural resolution limit of modularity.

A solution to bypass this problem is to adopt a hierarchical approach, by repeating for each community obtained separately, the whole procedure. The question then is to decide where to stop the recursion. We illustrate this issue with a variant of model (11) involving nested communities:

$$X_i(t) = a \alpha(t) + b_i \beta_{\sigma_i}(t) + b'_i \beta'_{\sigma'_i}(t) + c \gamma_i(t) \quad (12)$$

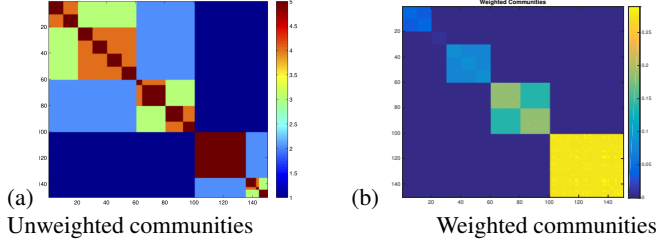
where the group dependent term now depends on two levels of embedded communities,  $\sigma_i$  and  $\sigma'_i$ . The results of a hierarchical approach yield the matrix of communities of Fig. 5 where each coefficient  $\text{Comm}_{kl}$  codes for the number of iterations up to which nodes  $k$  and  $l$  are kept in the same community. Although it allows to identify small communities, including the embedded ones, a problem is that there is no stopping criterion. For instance, the fourth community (nodes 61 to 80) is erroneously split into 4 sub-communities instead of just two if the iteration process were to stop at the third hierarchical level. Similarly, without control, this hierarchical approach forces the fifth monolithic community to fragment into meaningless groups. Conversely though, it is necessary to iterate the process up to level four to identify the small (sub-)communities formed by nodes 1 to 60.

## 4. THE PROPOSED ALGORITHM

In this section, we present our contributions to overcome the two issues of correlation based community detection, described above.

### 4.1. Mitigate the effect of global mode

A straightforward solution to limit the effect of global mode, as regards of a single dominant large community and of bulk squeezing,



**Fig. 6.** Comparison of the two possible outputs of **Comm**, the unweighted one and the weighted one from eq. (14), for a model correlation matrix identical to Fig. 5. The weighted **Comm** is similar to the expected correlation matrix, and the nested communities.

is to remove it before estimating the correlation matrix. The proposed approach is to first detect communities from the unmodified  $\mathbf{C}^{(s)}$ , using the Louvain algorithm. If only one community emerges, we remove the average behaviour from the time-series:

$$X_i(t) \leftarrow X_i(t) - X^{(g)}(t) \text{ where } X^{(g)}(t) = \frac{1}{N_A} \sum_{k \in A} X_k(t). \quad (13)$$

$\mathcal{V}$  (resp.  $N_{\mathcal{V}}$ ) stands for the set (resp. the cardinal) of nodes to be considered: initially all, and only a subset of them at subsequent iterations of the hierarchical algorithm. A new correlation matrix is then computed from the detrended  $N_{\mathcal{V}}$  time series. Without the global mode, the bulk is less squeezed towards the origin and the bounds  $\lambda_{\pm}$  are more accurately estimated. This is clearly illustrated on Fig. 3 (b), where the distribution of the eigenvalues of  $\mathbf{C}$  is displayed after the global mode was removed using Eq. (13), and to be compared with plot (a) showing the distribution before filtering.

#### 4.2. A weighted hierarchical approach

As described in Section 3.2, at each step of the hierarchical procedure, the community matrix entry  $\mathbf{Comm}_{kl}$  is incremented by one whenever the nodes  $k$  and  $l$  remain in the same cluster. Here, for each iteration on a newly found community  $\mathcal{V}$ , we propose to increment the matrix entries corresponding to any two nodes kept together in  $\mathcal{V}$ , by the modularity  $Q_{C_{\sigma_{\mathcal{V}}}}$ , associated to the new (embedded) detected partition of  $\sigma_{\mathcal{V}}$ :

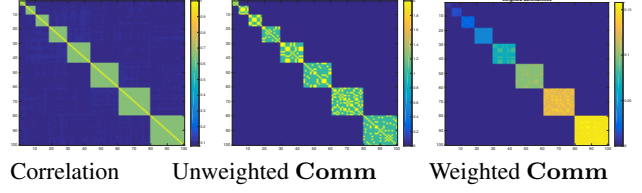
$$\mathbf{Comm}_{kl} \leftarrow \mathbf{Comm}_{kl} + Q_{C_{\sigma_{\mathcal{V}}}}. \quad (14)$$

For clusters that do not subdivide at finer resolution, the value  $Q_{C_{\sigma_{\mathcal{V}}}}$  will remain very small and the matrix entries  $\mathbf{Comm}_{kl}$  will stabilise. On the other hand, if a new subdivision emerges as the resolution increases, the matrix entries will undergo a larger increment and the corresponding sub-community will clearly distinguish in the community matrix **Comm**. Finally, the output is a weighted representation of hierarchical communities, as illustrated in Fig. 6.

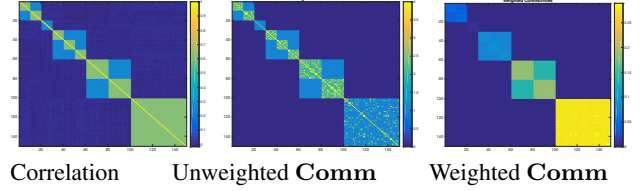
#### 4.3. Algorithm

Putting together the improvements developed in the previous two sections, we propose an algorithm that leads to a multi-resolution and weighted communities detection from a correlation matrix. The pseudo-code is as follows:

Input  $\{X_i(t), i = 1, \dots, N; t = 1, \dots, T\}$   
 $H$  (maximal depth of the hierarchy).  
Set  $h = 1$  and  $\mathcal{V} = \{1, \dots, N\}$ .



**Fig. 7.** Output (Weighted or not) of the algorithm for  $H = 2$  and comparison to the model correlation on the left.



**Fig. 8.** Output (Weighted or not) of the algorithm for  $H = 4$  in a case of nested communities, and comparison to the model correlation.

- (1) Consider the group of nodes in  $\mathcal{V}$ , and note  $N_{\mathcal{V}} = |\mathcal{V}|$ .
- (2.a) Compute and diagonalize  $\mathbf{C}$ .
- (2.b) From MP, filter out the eigen-components larger than  $\lambda_+ = (1 + \sqrt{N_{\mathcal{V}}/T})^2$  and form  $\mathbf{C}^{(s)}$ .
- (3.a) Apply Louvain algorithm on  $\mathbf{C}^{(s)}$  to find  $\sigma$ .
- (3.b) If there is only 1 community, remove the global mode as in Eq. (13).  
Go back to step (2.a) (or stop if all  $X_k$  are 0).
- (4) For each detected community  $\mathcal{V}^{\kappa}$ , increment the community matrix according to:

$$\forall k \in \mathcal{V}^{\kappa} \forall l \in \mathcal{V}^{\kappa} \quad \mathbf{Comm}_{kl} = \mathbf{Comm}_{kl} + Q_{C_{\mathcal{V}^{\kappa}}}(\sigma)$$

- (5) For each community  $\mathcal{V}^{\kappa}$ , set  $h = h + 1$  and repeat the procedure (1)-(5) for  $\mathcal{V} = \mathcal{V}^{\kappa}$ , until  $h = H$ .

## 5. EXAMPLES

### 5.1. Simple simulated examples

We illustrate the result of the previous algorithm on some simple examples. The first one follows the model of eq. (11), with heterogeneous communities in sizes, as in Fig. 4. The result is shown on Fig. 7 and it appears that the communities are all perfectly recovered.

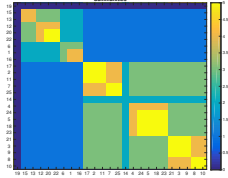
Then, the situation of Fig. 6 where communities are embedded in a hierarchical way as per eq. (12), is explored. The result is shown on Fig. 8. Here again, the weighted approach with our algorithm outputs a correct multi resolution representation of the communities.

### 5.2. Example with real data

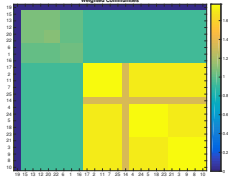
We consider now data from the environmental sensors from the Live E! project [11] for illustration. The data consists in several time series of temperature, with a time resolution set to 10 minutes. We will explore a specific zone, the 25 sensors in Kurashiki city, Okayama Prefecture, Japan. The goal is to group together sensors experiencing similar temperature evolutions. As expected (and seen on Fig. 1), there is a dominating global mode for all the sensors and it is the fluctuations around this mode that are interesting. That justifies the use



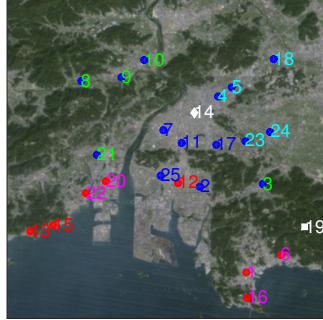
Unweighted communities



Weighted communities



Map with sensors



**Fig. 9.** Output of the algorithm for the Live E! temperature sensors. The communities displayed on the left are used to color the symbol of the sensor position according to the first level, and to color the number of the sensor according to the second level of the hierarchy of communities. The two sensors in white are malfunctioning sensors and detected as such outside communities.

of the approach developed here. Figure 9 shows both the output with unweighted and weighted community matrix.

In both cases, one sensor (19) is always separated from the others: in fact this is an indication of the malfunctions at that time of this sensor, and the same can be said for sensor (14) which is an outlier also at the following levels. Other sensors are then grouped in two large communities, one with 8 sensors which are relatively homogeneous, and the second with the remaining 15 sensors appear to have 3 sub-communities (plus 14 as outlier).

The finding is that the communities are essentially geographical in their positioning: the separation in 2 big communities separate the sensors in places near the sea from sensors more inland, and the following levels are associated to refinements depending whether there is a part of forest near the sensor's location or not. This is relevant as sea and forests have a major impact on temperature, with for instance a cooling effect that reduces the possible fluctuations in temperature measured by these sensors around the global mode.

## 6. DEVELOPMENTS AND PERSPECTIVES

The proposed algorithm relies on the simplification of the correlation matrix via its decomposition and the removal of its random part. In the spectral domain, the action is to filter the matrix by keeping only its largest eigenvalues. The direct implementation of that is to diagonalize  $\mathbf{C}$ , and this could be cumbersome for large problems. An alternative is to realize that keeping only  $\mathbf{C}^{(s)}$  is a low-pass filtering<sup>1</sup> of  $\mathbf{C}$ , that keeps only eigenvalues larger than  $\lambda_+$ . This interpretation in terms of filtering is possible because  $\mathbf{C}$  is definite positive: it is diagonalizable with real positive eigenvalues; hence, contrary to graph filters for adjacency matrix, as proposed in [14], its spectrum is well defined and ordered. Instead of computing  $\mathbf{C}$ , one can try to compute the effect of this filtering in the spectral domain (using the method of [15]) of  $\mathbf{C}$  (as done in [7] for community detection with wavelets) and estimate  $\mathbf{C}^{(s)}$  by applying this filtering to some

fixed vectors. We do not detail this graph filter implementation of the method further on here, as it remains a work in progress.

Perspectives of the present work of community mining in correlation matrices would be first to go to problems of larger scales (as the MP distribution would remain well valid), using the mentioned graph filter implementation, and, second, to be able to take also into account some adjacency matrix in the study, e.g., position in spaces (and nearest neighbors) for sensor networks. This will be studied in future works.

## 7. REFERENCES

- [1] S. Fortunato, "Community detection in graphs," *Physics Reports*, vol. 486, no. 3-5, pp. 75–174, 2010.
- [2] M. Newman, "Modularity and community structure in networks," *PNAS*, vol. 103, no. 23, p. 8577, 2006.
- [3] M. MacMahon and D. Garlaschelli, "Community detection for correlation matrices," *Phys. Rev. X*, vol. 5, p. 021005, 2015.
- [4] S. Fortunato and M. Barthélemy, "Resolution limit in community detection," *PNAS*, vol. 104, no. 1, p. 36, 2007.
- [5] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [6] M. Rosvall and C. Bergstrom, "Mapping change in large networks," *PloS one*, vol. 5, no. 1, p. e8694, 2010.
- [7] N. Tremblay and P. Borgnat, "Graph wavelets for multiscale community mining," *Signal Processing, IEEE Transactions on*, vol. 62, no. 20, pp. 5227–5239, Oct 2014.
- [8] R. R. Nadakuditi and M. E. J. Newmann, "Graph spectra and the detectability of community structure in networks," *Phys. Rev. Lett.*, vol. 108, p. 188701, 2012.
- [9] H. T. Ali and R. Couillet, "Performance analysis of spectral community detection in realistic graph models," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP'16)*, 2016, (submitted to).
- [10] V. Blondel, J. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *JSTAT*, vol. 2008, no. 10, p. P10008, 2008.
- [11] S. Matsuura, H. Ishizuka, H. Ochiai, S. Doi, S. Ishida, M. Nakayama, H. Esaki, and H. Sunahara, "Live E! project: Establishment of infrastructure sharing environmental information," in *SAINT Workshops*, Jan. 2007, pp. pp.67–67.
- [12] J. P. Bouchaud and M. Potters, *Theory of Financial Risk and Derivative Pricing*, 2nd ed. CUP, 2003.
- [13] R. Couillet and M. Debbah, *Random Matrix Methods for Wireless Communications*. CUP, 2011.
- [14] A. Sandryhaila and J. Moura, "Discrete signal processing on graphs," *Signal Processing, IEEE Transactions on*, vol. 61, no. 7, pp. 1644–1656, April 2013.
- [15] D. Hammond, P. Vandergheynst, and R. Gribonval, "Wavelets on graphs via spectral graph theory," *Applied and Computational Harmonic Analysis*, vol. 30, no. 2, pp. 129–150, 2011.

<sup>1</sup>The largest eigenvalues are associated to the more global mode, hence the equivalent "frequency" is ordered as the opposite of the  $\lambda_k$  and low "frequency" are for large eigenvalues.